# Software work in mathematical biology at Georgia Tech

Maxie Dion Schmidt

maxieds@gmail.com
http://people.math.gatech.edu/~mschmidt34
https://github.com/maxieds

Sandia National Labs
Technical Presentation
Spring 2022

# Introduction – Applications of RNA research

- ▶ RNA sequencing utilized in cancer research and therapy
- ▶ mRNA vaccines are newly available to the public but have been studied for decades
- ▶ Allowed for the rapid development of a COVID-19 vaccine
- ▶ Research on RNA may eventually play a central role in medical applications

# What is mathematical biology?

▶ Mathematical biology (MathBio) uses mathematical models as theoretical abstractions of the natural structure of living organisms

▶ In this talk we will discuss my work as a software engineering RA with the *Georgia Tech Discrete Mathematics and Molecular Biology* group (gtDMMB)

# RNA basics

# RNA basics

- RNA is a single-stranded molecule similar to DNA
- A strand of RNA has a backbone of alternating sugar (ribose) and phosphate groups
- Each sugar has one of four base types attached to it (**A**–**U**–**C**–**G**)
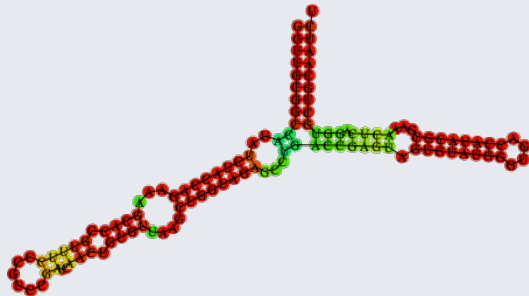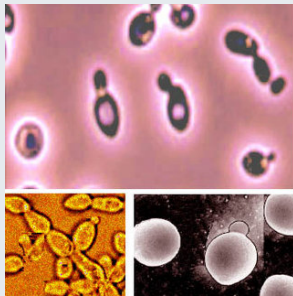- Each of the **A**–**U**–**C**–**G** bases can fold to form bonds in pairs

# Arc diagrams – Discussion example – S. Cerevisiae (yeast)

GGUUGCGGCCAUAUCUACCAGAAAGCACCGUUUCCCGUCCGAUCAACUGUGUUAAGCUGGUAGA

(((((((((....(((((((((...((((.((((((......))..))).))))....)))))))))

GCCUGACCGAGUAGUGUAUGGGUGACCAUACGCGAAACUCAGGUGCUGCAAUCU

.....(((((((.((((((((....))))))))...)))).))))))))).



**(Actual microscopic views)**     **(Radial view of 2D MFE structure)**

# RNA secondary structures

▶ RNA base sequence and can have more than one 2D or 3D structure

▶ Obtaining the 1D structure for organisms is easy via modern sequencing

▶ Characterizing 3D molecular conformations is still comparatively hard

▶ Understanding the 2D secondary structures (base pairings) remains a crucial component of ribonomics research

▶ RNA folding prediction programs generate 2D structures given the 1D base sequence
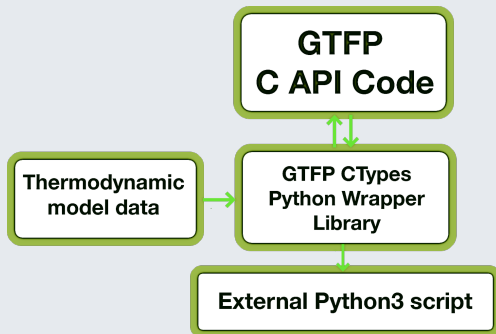
# GTFoldPython software project

# GTFold – Overview I

▶ Accurate and efficient RNA secondary structure prediction is an important open problem in computational molecular biology

▶ **Input:** The base sequence (FASTA) of the organism
**Output:** MFE or suboptimal secondary structures (base pair data)

▶ GTFold is the first implementation of RNA secondary structure prediction by thermodynamic optimization for modern multi-core computers

▶ The speedup is useful to researchers working with very RNA sequences, such as RNA viral genomes

# GTFold – Overview II

- ▶ Original GTFold produced command line only utilities
- ▶ Motivated by the need for a Python interface to get the same GTFold data
- ▶ Started work in the Fall of 2019 writing Python3 bindings for GTFold
- ▶ Hardest part of development work? Requirements for highly robust custom cross platform build scripts

# GTFoldPython – Introduction

▶ GTFoldPython (GTFP): Python3 bindings based around the original GTFold sources written in C++

▶ Backend uses the Python3 C API

▶ Frontend interface is a wrapper library that uses `CTypes` to call the C API functions

# GTFoldPython – Comparison – Python C API code

```
1   PyObject * GetMFEStructure(const char *baseSeq, ConsListCType_t consList, int consLength) {
2       /* Error checking omitted ... */
3       MFEStructRuntimeArgs_t rtArgs;
4       InitMFEStructRuntimeArgs(&rtArgs);
5       rtArgs.baseSeq = baseSeq;
6       SetRTArgsSequenceLength(rtArgs, strlen(baseSeq));
7       if(ParseGetMFEStructureArgs(consList, consLength, &rtArgs) != GTFPYTHON_ERRNO_OK) {
8           FreeMFEStructRuntimeArgs(&rtArgs);
9           return ReturnPythonNone();
10      }
11      if(InitGTFoldMFEStructureData(&rtArgs) != GTFPYTHON_ERRNO_OK) {
12          FreeMFEStructRuntimeArgs(&rtArgs);
13          return ReturnPythonNone();
14      }
15      double mfe = ComputeMFEStructure(&rtArgs);
16      if(GetLastErrorCode() != GTFPYTHON_ERRNO_OK) {
17          return ReturnPythonNone();
18      }
19      if(WRITEAUXFILES) {
20          ConfigureOutputFileSettings();
21          save_ct_file(outputFile, baseSeq, mfe);
22      }
23      char *dbMFEStruct = ComputeDOTStructureResult(rtArgs.numBases);
24      PyObject *mfeTupleRes = PrepareMFETupleResult(mfe, dbMFEStruct);
25      Free(dbMFEStruct);
26      FreeMFEStructRuntimeArgs(&rtArgs);
27      FreeGTFoldMFEStructureData(rtArgs.numBases);
28      if(mfeTupleRes == NULL) {
29          return ReturnPythonNone();
30      }
31      return mfeTupleRes;
```

# GTFoldPython – Comparison – Wrapper library code

```python
1    ## Library initialization code:
2    if GTFPConfig.PLATFORM_DARWIN:
3        GTFoldPython._libGTFoldHandle = ctypes.cdll.LoadLibrary("GTFoldPython.dylib")
4    else:
5        GTFoldPython._libGTFoldHandle = ctypes.PyDLL("GTFoldPython.so",
6                                          mode=ctypes.RTLD_GLOBAL, use_errno=True)
7    @staticmethod
8    def _WrapCTypesFunction(funcname, restype=None, argtypes=None):
9        return GTFoldPython._libGTFoldHandle.__getattr__(funcname)
10
11   @staticmethod
12   def GetMFEStructure(baseSeq, consList = []):
13       """Get the MFE and MFE structure (in DOTBracket structure notation)
14       :param baseSeq: A string of valid bases (ATGU/X)
15       :param consList: A list of constraints on the MFE structure
16       :return: A tuple (MFE as double, MFE structure as string in DOTBracket notation)
17       :rtype: tuple
18       """
19       GTFoldPython._ConstructLibGTFold()
20       resType = ctypes.py_object
21       argTypes = [ GTFPTypes.CStringType,
22                    GTFPTypes.FPConstraintsListType(consList),
23                    ctypes.c_int ]
24       libGTFoldFunc = GTFoldPython._WrapCTypesFunction("GetMFEStructure", resType, argTypes)
25       (mfe, mfeStruct) = libGTFoldFunc(GTFPTypes.CString(baseSeq),
26                                        GTFPTypes.FPConstraintsList(consList),
27                                        len(consList))
28       return (float(mfe), str(mfeStruct))
```

# GTFoldPython – Example – Find MFE and MFE structure

**External Python3 script source:**

```python
1   import sys, os
2   from GTFoldPythonImportAll import *
3
4   GTFP.Init()
5   GTFP.Config(quiet = False, debugging = False, verbose = False, stdmsgout = "stderr")
6
7   baseSeqFPCons = "GCAUUGGAGAUGGCAUUCCUCCAUUAACAAACCGCUGCGCCCGUAGCAGCUGAUGAUGCCUACAGA"
8   consListFP = GTFPUtils.ReadFPConstraintsFromFile("../Testing/TestData/tRNA/yeast.fa.cons")
9
10  (mfe, mfeDOTStruct) = GTFP.GetMFEStructure(baseSeqFPCons, consListFP)
11  print("MFE_%1.3f_=>_MFE_DOT_STRUCT_\"%s\"\n\n" % (mfe, mfeDOTStruct))
```

**Terminal output printed upon invoking the script above:**

```
1   MFE -17.200 => MFE DOT STRUCT "(((((((((.........))))...........((((.......)))))....)))))......
```

# The RNAStructViz application

# RNAStructViz: Graphical base pairing analysis

▶ RNAStructViz was a project developed by Professor Christine Heitsch and Dr. S. Cheney

▶ My first project working with the gtDMMB group in the Summer of 2018
  - Modernize the C++ source
  - Add support for enhanced graphics using the `cairo` library
  - Re-write the dated build scripts
  - Improve and support the project in the long term

▶ Key feature of RNAStructViz: Visualization and comparisons via arc diagrams of RNA secondary structures

# RNAStructViz – Comparison of features

| ↓ Feature sets | RNAStructViz | FORNA | jViz.RNA | R-chie | RNAbows | VARNA |
|---|---|---|---|---|---|---|
| **Software support →** | | | | | | |
| **❶ — Platform and availability —** | | | | | | |
| Mac OSX support | ✔ | ✓ | ✓ | ✓ | ✓ | ✓ |
| Linux / Unix support | ✔ | ✓ | ✓ | ✓ | ✓ | ✓ |
| Windows support | ✗ | ✓ | ✓ | ✓ | ✓ | ✓ |
| Open source software | ✔ | ✓ | ✓ | ✓ | ✗ | ✓ |
| Requires external libraries | ✔ | ✓ | ✓ | ✓ | ✗ | ✓ |
| **❷ — Software usability criteria —** | | | | | | |
| Graphical user interface | ✔ | ✓ | ✓ | ✓ | ✓ | ✓ |
| Web interface | ✗ | ✓ | ✗ | ✓ | ✓ | ✗ |
| Multi-window interface | ✔ | ✗ | ✗ | ✗ | ✗ | ✗ |
| Compares 2 structures at once | ✔ | ✓ | ✗ | ✓ | ✓ | ✗ |
| Compares 3 structures at once | ✔ | ✓* | ✗ | ✗ | ✗ | ✗ |

| ↓ Feature sets | RNAStructViz | FORNA | jViz.RNA | R-chie | RNAbows | VARNA |
|---|---|---|---|---|---|---|
| **Software support →** | | | | | | |
| **❸ — Support for standard formats —** | | | | | | |
| CT files | ✔ | ✗ | ✓ | ✓ | ✗ | ✓ |
| Dot-bracket files | ✔ | ✗ | ✓ | ✓ | ✗ | ✓ |
| Built-in file viewer | ✔ | ✗ | ✗ | ✗ | ✗ | ✓ |
| Requires specialized format | ✗ | ✓ | ✗ | ✗ | ✗ | ✗ |
| Can edit sequence data | ✗ | ✓ | ✓ | ✗ | ✓ | ✓ |
| **❹ — Views and diagram type support —** | | | | | | |
| Has comparison statistics | ✔ | ✗ | ✓** | ✓ | ✓ | ✗ |
| Plots circular arc diagrams | ✔ | ✗ | ✓ | ✓ | ✓ | ✗ |
| Plots radial diagrams | ✔ | ✓ | ✓ | ✗ | ✗ | ✓ |

*A comparison of selected features across related tools; an extended survey appears in the RNAStructViz WIKI.*

# RNAStructViz Screenshot – Loading sample structures I

# RNAStructViz Screenshot – Loading sample structures II
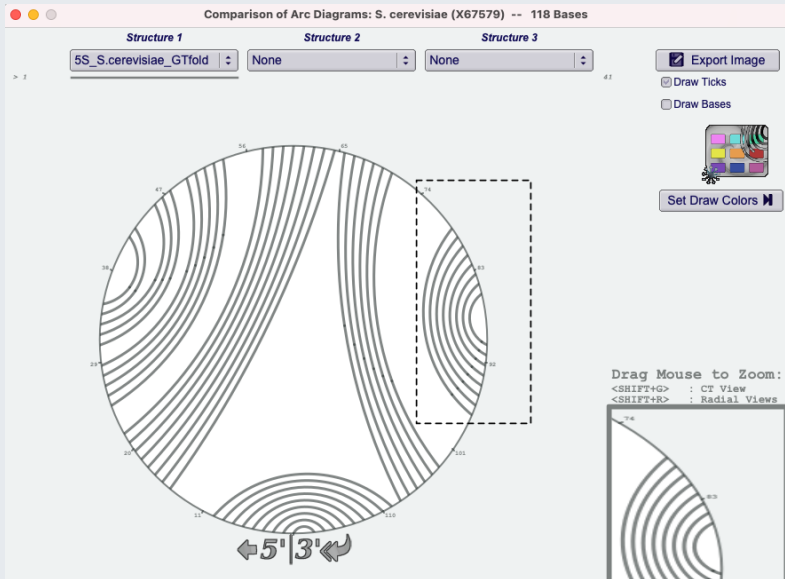
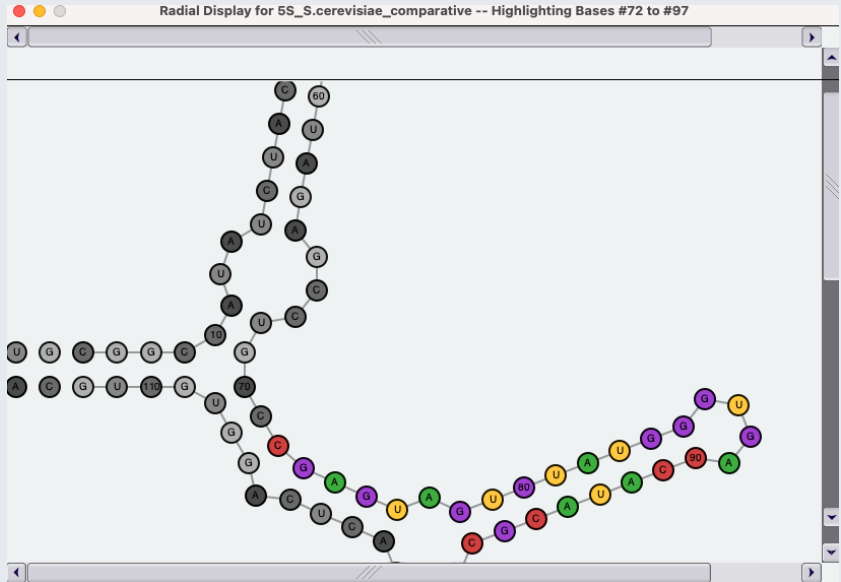# RNAStructViz Screenshot – Arc diagram window

# Arc diagram window – Discussion

▶ The bases from position #1 to #LengthOfBaseSequenceString at equidistant spacings around a circle

▶ The sequentially numbered base pairs are ordered around the circle counter-clockwise starting from the bottom

▶ An arc connecting paired bases is drawn within the circle

# Arc diagram window – Zoom select

# Arc diagram zoom – Radial layout visualization



Radial Display for 5S_S.cerevisiae_comparative -- Highlighting Bases #72 to #97

# Arc diagram zoom – CT segment visualization

# Arc diagram window – Comparing multiple structures

# RNAStructViz Screenshot – Statistics window

# Wrapping up

# Summary of accomplishments with gtDMMB software I

- ▶ Success modernizing and enhancing the source code for these projects in computational and mathematical biology
- ▶ Success in modernizing and extending build scripts to support installation on MacOS, Linux and Unix-based systems
- ▶ A few of the software projects we worked on:
  - RNAStructViz
  - GTFold (CMake for MacOS and Linux)
  - GTFoldPython

# Summary of accomplishments with gtDMMB software II

▶ Application note re-introducing our new work on RNAStructViz published in *Bioinformatics* in 2021

▶ Sister RNA labs that helped with testing and/or use our software include:

  - Computational RNA Genomics Lab at University of California Davis
  - Laederach Lab at the University of North Carolina at Chapel Hill
  - Mathews Lab at the University of Rochester

# Concluding remarks

# The End

Questions?

Comments?

Feedback?

# Thank you for your time!